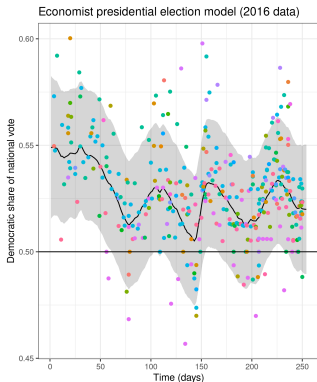


Approximate data deletion and replication with the Bayesian influence function

Ryan Giordano (rgiordano@berkeley.edu, UC Berkeley), Tamara Broderick (MIT)
Stanford Statistics Seminar May 2024

Economist 2016 Election Model [Gelman and Heidemanns, 2020]



A time series model to predict the 2016 US presidential election outcome from polling data.

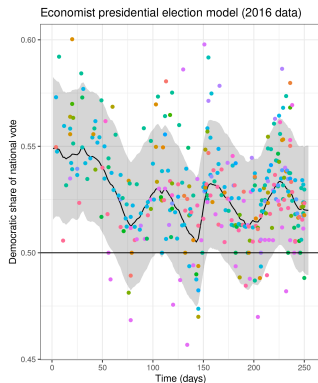
Model:

- $X = x_1, \dots, x_N =$ Polling data ($N = 361$).
- $\theta =$ Lots of random effects (day, pollster, etc.)
- $f(\theta) =$ Democratic % of vote on election day

Typically, we compute Markov chain Monte Carlo (MCMC) draws from the posterior $p(\theta|X)$.

We want to know $\mathbb{E}_{p(\theta|X)} [f(\theta)]$.

Economist 2016 Election Model [Gelman and Heidemanns, 2020]



A time series model to predict the 2016 US presidential election outcome from polling data.

Model:

- $X = x_1, \dots, x_N =$ Polling data ($N = 361$).
- $\theta =$ Lots of random effects (day, pollster, etc.)
- $f(\theta) =$ Democratic % of vote on election day

Typically, we compute Markov chain Monte Carlo (MCMC) draws from the posterior $p(\theta|X)$.

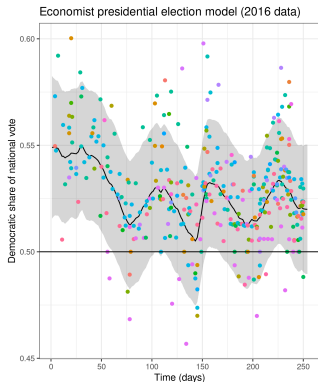
We want to know $\mathbb{E}_{p(\theta|X)} [f(\theta)]$.

The people who responded to the polls were randomly selected.

If we had selected a different random sample, how much would our estimate have changed?

Idea: Re-fit with bootstrap samples of data [Huggins and Miller, 2023]

Economist 2016 Election Model [Gelman and Heidemanns, 2020]



A time series model to predict the 2016 US presidential election outcome from polling data.

Model:

- $X = x_1, \dots, x_N =$ Polling data ($N = 361$).
- $\theta =$ Lots of random effects (day, pollster, etc.)
- $f(\theta) =$ Democratic % of vote on election day

Typically, we compute Markov chain Monte Carlo (MCMC) draws from the posterior $p(\theta|X)$.

We want to know $\mathbb{E}_{p(\theta|X)} [f(\theta)]$.

The people who responded to the polls were randomly selected.

If we had selected a different random sample, how much would our estimate have changed?

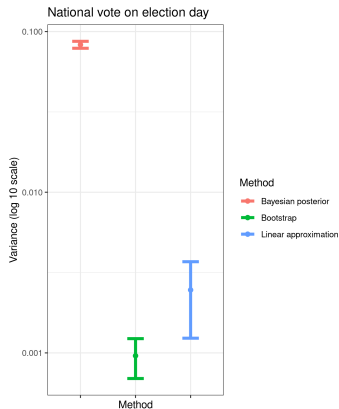
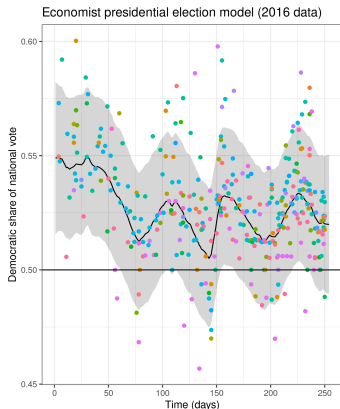
Idea: Re-fit with bootstrap samples of data [Huggins and Miller, 2023]

Problem: Each MCMC run takes about 10 hours (Stan, six cores).

Proposal: Use full-data posterior draws to form a linear approximation to *data reweightings*.

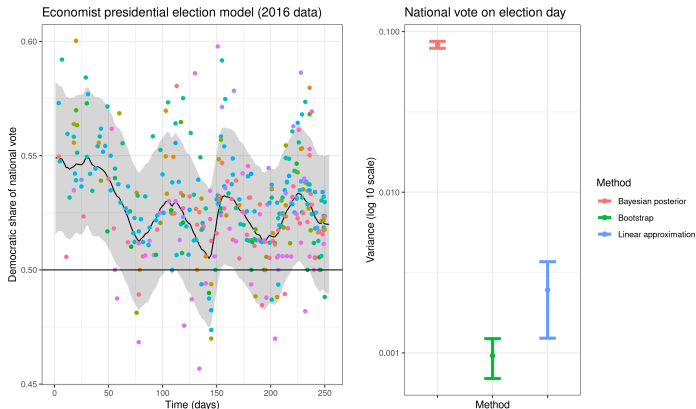
Results

Proposal: Use full-data posterior draws to form a linear approximation to *data reweightings*.



Results

Proposal: Use full-data posterior draws to form a linear approximation to *data reweightings*.



Compute time for 100 bootstraps: 51 days

Compute time for the linear approximation: Seconds
(But note the approximation has some error)

- Data reweighting
 - Write the change in the posterior expectation as **linear component** + **error**
 - The **linear component** can be computed from a single run of MCMC

- Data reweighting
 - Write the change in the posterior expectation as **linear component** + **error**
 - The **linear component** can be computed from a single run of MCMC
- Finite-dimensional problems with posteriors which concentrate asymptotically
 - As $N \rightarrow \infty$, the linear component provides an arbitrarily good approximation

- Data reweighting
 - Write the change in the posterior expectation as **linear component** + **error**
 - The **linear component** can be computed from a single run of MCMC
- Finite-dimensional problems with posteriors which concentrate asymptotically
 - As $N \rightarrow \infty$, the linear component provides an arbitrarily good approximation
- High-dimensional problems
 - The linear component is the same order as the error
 - Even for parameters which concentrate, even as $N \rightarrow \infty$

- Data reweighting
 - Write the change in the posterior expectation as **linear component** + **error**
 - The **linear component** can be computed from a single run of MCMC
- Finite-dimensional problems with posteriors which concentrate asymptotically
 - As $N \rightarrow \infty$, the linear component provides an arbitrarily good approximation
- High-dimensional problems
 - The linear component is the same order as the error
 - Even for parameters which concentrate, even as $N \rightarrow \infty$
- What should the exchangeable unit be?

Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

Original weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

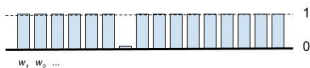
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

Original weights:



Leave-one-out weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

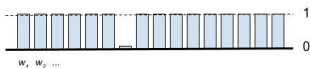
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

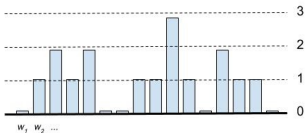
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

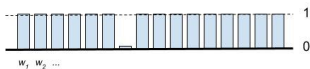
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

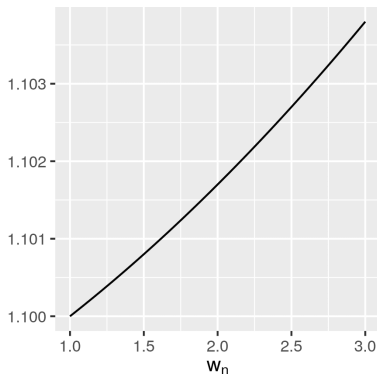
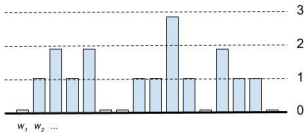
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

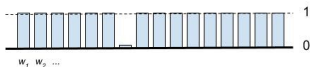
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

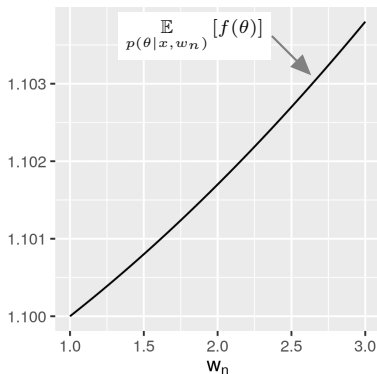
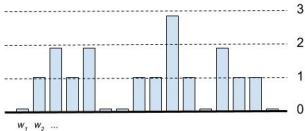
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

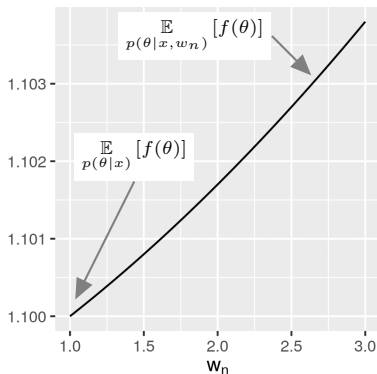
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

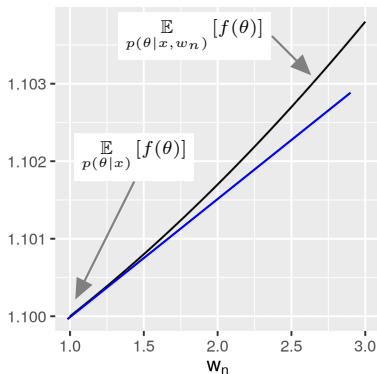
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

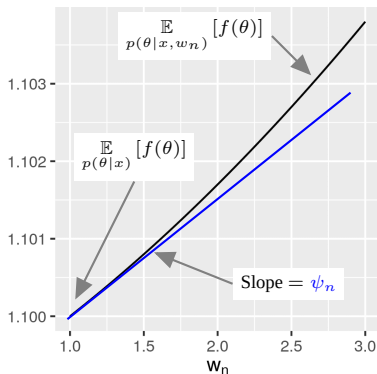
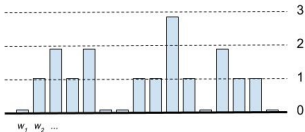
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

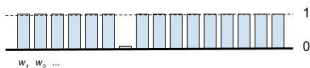
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

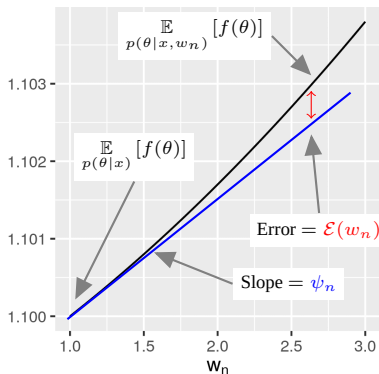
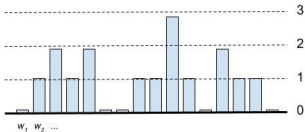
Original weights:



Leave-one-out weights:



Bootstrap weights:



Data re-weighting.

Augment the problem with *data weights* w_1, \dots, w_N . We can write $\mathbb{E}_{p(\theta|X,w)}[f(\theta)]$.

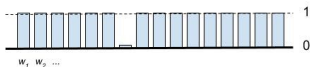
$$\ell_n(\theta) := \log p(x_n|\theta)$$

$$\log p(X|\theta, w) = \sum_{n=1}^N w_n \ell_n(\theta)$$

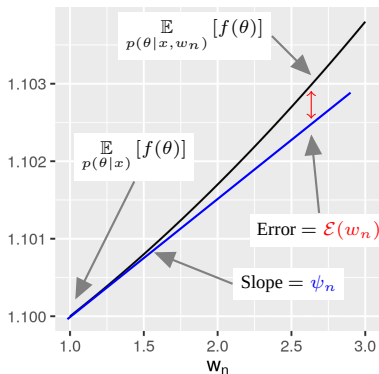
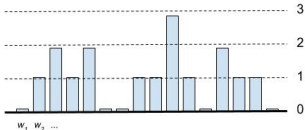
Original weights:



Leave-one-out weights:



Bootstrap weights:



The re-scaled slope $N\psi_n$ is known as the “influence function” at data point x_n .

$$\mathbb{E}_{p(\theta|X,w)}[f(\theta)] - \mathbb{E}_{p(\theta|X)}[f(\theta)] = \sum_{n=1}^N \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

How can we use the approximation?

Assume the **slope** is computable and **error** is small.

$$\mathbb{E}_{p(\theta|X,w)}[f(\theta)] - \mathbb{E}_{p(\theta|X)}[f(\theta)] = \sum_{n=1}^N \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

How can we use the approximation?

Assume the **slope** is computable and **error** is small.

$$\mathbb{E}_{p(\theta|X,w)}[f(\theta)] - \mathbb{E}_{p(\theta|X)}[f(\theta)] = \sum_{n=1}^N \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

Bootstrap. Draw bootstrap weights $w \sim p(w) = \text{Multinomial}(N, N^{-1})$.

$$\begin{aligned} \text{Bootstrap variance} &= \text{Var}_{p(w)} \left(\mathbb{E}_{p(\theta|X,w)}[f(\theta)] \right) \\ &= \text{Var}_{p(w)} \left(\sum_{n=1}^N \psi_n(w_n - 1) + \mathcal{E}(w_n) \right) \\ &= \frac{1}{N^2} \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2 + \text{Term involving } \mathcal{E}(w_n) \text{ for } n = 1, \dots, N \\ &\approx \frac{1}{N^2} \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2 \end{aligned}$$

Expressions for the slope and error

How to compute the slopes ψ_n ? How large is the error $\mathcal{E}(w)$?

For simplicity, let us consider a single weight for the moment.

$$\mathbb{E}_{p(\theta|X, w_n)} [f(\theta)] - \mathbb{E}_{p(\theta|X)} [f(\theta)] = \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

Expressions for the slope and error

How to compute the slopes ψ_n ? How large is the error $\mathcal{E}(w)$?

For simplicity, let us consider a single weight for the moment.

$$\mathbb{E}_{p(\theta|X, w_n)} [f(\theta)] - \mathbb{E}_{p(\theta|X)} [f(\theta)] = \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

Let an overbar denote “posterior–mean zero.” For example, $\bar{f}(\theta) := f(\theta) - \mathbb{E}_{p(\theta|X)} [f(\theta)]$.

By dominated convergence and the mean value theorem, for some $\tilde{w}_n \in [0, w_n]$:

$$\psi_n = \underbrace{\mathbb{E}_{p(\theta|X)} [\bar{f}(\theta) \bar{\ell}_n(\theta)]}_{\text{Estimatable with MCMC!}} \quad \mathcal{E}(w_n) = \frac{1}{2} \underbrace{\mathbb{E}_{p(\theta|X, \tilde{w}_n)} [\bar{f}(\theta) \bar{\ell}_n(\theta) \bar{\ell}_n(\theta)]}_{\text{Cannot compute directly (don't know } \tilde{w})} (w_n - 1)^2$$

Expressions for the slope and error

How to compute the slopes ψ_n ? How large is the error $\mathcal{E}(w)$?

For simplicity, let us consider a single weight for the moment.

$$\mathbb{E}_{p(\theta|X, w_n)} [f(\theta)] - \mathbb{E}_{p(\theta|X)} [f(\theta)] = \psi_n(w_n - 1) + \mathcal{E}(w_n)$$

Let an overbar denote “posterior–mean zero.” For example, $\bar{f}(\theta) := f(\theta) - \mathbb{E}_{p(\theta|X)} [f(\theta)]$.

By dominated convergence and the mean value theorem, for some $\tilde{w}_n \in [0, w_n]$:

$$\begin{aligned} \psi_n &= \underbrace{\mathbb{E}_{p(\theta|X)} [\bar{f}(\theta) \bar{\ell}_n(\theta)]}_{\text{Estimatable with MCMC!}} & \mathcal{E}(w_n) &= \frac{1}{2} \underbrace{\mathbb{E}_{p(\theta|X, \tilde{w}_n)} [\bar{f}(\theta) \bar{\ell}_n(\theta) \bar{\ell}_n(\theta)]}_{\text{Cannot compute directly (don't know } \tilde{w})} (w_n - 1)^2 \\ &= O_p(N^{-1}) \text{ under posterior concentration} & &= O_p(N^{-2}) \text{ under posterior concentration} \end{aligned}$$

Expressions for the slope and error

How to compute the slopes ψ_n ? How large is the error $\mathcal{E}(w)$?

For simplicity, let us consider a single weight for the moment.

$$\mathbb{E}_{p(\theta|X, w_n)} [f(\theta)] - \mathbb{E}_{p(\theta|X)} [f(\theta)] = \psi_n (w_n - 1) + \mathcal{E}(w_n)$$

Let an overbar denote “posterior–mean zero.” For example, $\bar{f}(\theta) := f(\theta) - \mathbb{E}_{p(\theta|X)} [f(\theta)]$.

By dominated convergence and the mean value theorem, for some $\tilde{w}_n \in [0, w_n]$:

$$\begin{aligned} \psi_n &= \underbrace{\mathbb{E}_{p(\theta|X)} [\bar{f}(\theta) \bar{\ell}_n(\theta)]}_{\text{Estimatable with MCMC!}} & \mathcal{E}(w_n) &= \frac{1}{2} \underbrace{\mathbb{E}_{p(\theta|X, \tilde{w}_n)} [\bar{f}(\theta) \bar{\ell}_n(\theta) \bar{\ell}_n(\theta)]}_{\text{Cannot compute directly (don't know } \tilde{w})} (w_n - 1)^2 \\ &= O_p(N^{-1}) \text{ under posterior concentration} & &= O_p(N^{-2}) \text{ under posterior concentration} \end{aligned}$$

Theorem 1 [Giordano and Broderick, 2023] (paraphrase):

If the posterior $p(\theta|X)$ “concentrates” (e.g. as in the Bernstein–von Mises theorem),^a then

$$w_n \mapsto N \left(\mathbb{E}_{p(\theta|X, w_n)} [f(\theta)] - \mathbb{E}_{p(\theta|X)} [f(\theta)] \right)$$

becomes linear as $N \rightarrow \infty$, with slope $\lim_{N \rightarrow \infty} \psi_n$.

^aExisting results are sufficient for a *particular weight* [Kass et al., 1990]. Giordano and Broderick [2023] proves that the result holds when averaged over all weights, as needed for variance estimation.

Negative binomial experiment

Example: Negative binomial models with an unknown parameter γ .

For $n = 1, \dots, N$ let $x_n | \gamma \stackrel{iid}{\sim} \text{NegativeBinomial}(\alpha, \gamma)$ for fixed α .

$$\text{Write } \log p(X | \lambda, \gamma, w) = \sum_{n=1}^N w_n \ell_n(\gamma).$$

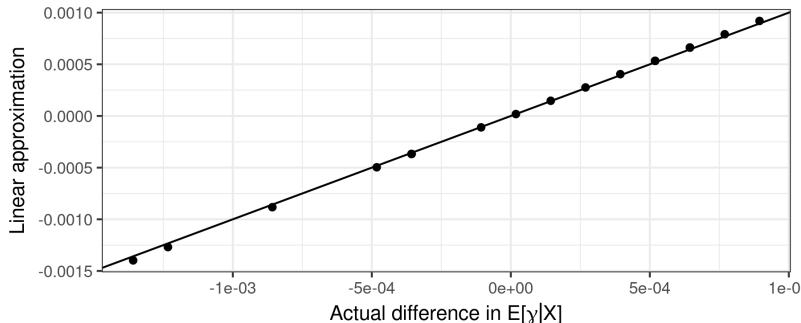
Negative binomial experiment

Example: Negative binomial models with an unknown parameter γ .

For $n = 1, \dots, N$ let $x_n | \gamma \stackrel{iid}{\sim} \text{NegativeBinomial}(\alpha, \gamma)$ for fixed α .

$$\text{Write } \log p(X | \lambda, \gamma, w) = \sum_{n=1}^N w_n \ell_n(\gamma).$$

Negative Binomial model
leaving out single datapoints with $N = 800$



Assumptions sketch:

- A well-behaved MAP *maximum a posteriori* estimator $\hat{\theta}$ exists:
 - The dimension of θ is fixed as $N \rightarrow \infty$.
 - The expected log likelihood has a unique maximum at θ_∞
 - The observed log likelihood satisfies $\hat{\theta} \rightarrow \theta_\infty$
 - The expected log likelihood Hessian \mathcal{I} is negative definite at θ_∞
- We can apply standard asymptotics:
 - The log prior and log likelihood are four times continuously differentiable
 - The prior is proper, and a technical set of squared expectations are finite
 - The log likelihood derivatives are dominated by a square-integrable envelope function in a neighborhood of θ_∞ .

Variance consistency theorem

Assumptions sketch:

- A well-behaved MAP *maximum a posteriori* estimator $\hat{\theta}$ exists:
 - The dimension of θ is fixed as $N \rightarrow \infty$.
 - The expected log likelihood has a unique maximum at θ_∞
 - The observed log likelihood satisfies $\hat{\theta} \rightarrow \theta_\infty$
 - The expected log likelihood Hessian \mathcal{I} is negative definite at θ_∞
- We can apply standard asymptotics:
 - The log prior and log likelihood are four times continuously differentiable
 - The prior is proper, and a technical set of squared expectations are finite
 - The log likelihood derivatives are dominated by a square-integrable envelope function in a neighborhood of θ_∞ .

Theorem 2 [Giordano and Broderick, 2023]:

Under the above assumptions,

$$\sqrt{N} \left(\mathbb{E}_{p(\theta|X)} [g(\theta)] - g(\theta_\infty) \right) \xrightarrow[N \rightarrow \infty]{dist} \mathcal{N}(0, V^g) \quad \text{and} \quad (1)$$
$$\frac{1}{N} \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2 \xrightarrow[N \rightarrow \infty]{prob} V^g.$$

Variance consistency theorem

Assumptions sketch:

- A well-behaved MAP *maximum a posteriori* estimator $\hat{\theta}$ exists:
 - The dimension of θ is fixed as $N \rightarrow \infty$.
 - The expected log likelihood has a unique maximum at θ_∞
 - The observed log likelihood satisfies $\hat{\theta} \rightarrow \theta_\infty$
 - The expected log likelihood Hessian \mathcal{I} is negative definite at θ_∞
- We can apply standard asymptotics:
 - The log prior and log likelihood are four times continuously differentiable
 - The prior is proper, and a technical set of squared expectations are finite
 - The log likelihood derivatives are dominated by a square-integrable envelope function in a neighborhood of θ_∞ .

Theorem 2 [Giordano and Broderick, 2023]:

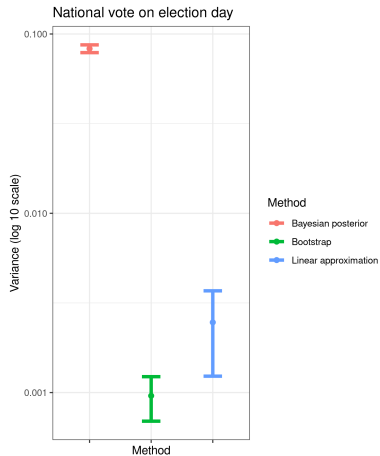
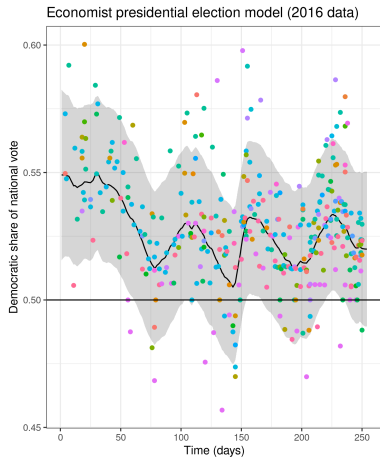
Under the above assumptions,

$$\sqrt{N} \left(\mathbb{E}_{p(\theta|X)} [g(\theta)] - g(\theta_\infty) \right) \xrightarrow[N \rightarrow \infty]{dist} \mathcal{N}(0, V^g) \quad \text{and} \quad (1)$$
$$\frac{1}{N} \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2 \xrightarrow[N \rightarrow \infty]{prob} V^g.$$

Equation 1 and the form of V^g is known ([Kleijn and Van der Vaart, 2012]).

Our contribution is a consistent estimator of V^g using posterior samples rather than $\hat{\theta}$.

How to connect to the election data?



Problem: MCMC is only interesting when the posterior doesn't concentrate.

Example: Exponential families with random effects (REs) λ and fixed effects γ .

Example: Exponential families with random effects (REs) λ and fixed effects γ .

If the observations per random effect remains bounded as $N \rightarrow \infty$, then

- Parameter λ (“local”) grows in dimension with N .
- Parameter γ (“global”) is finite-dimensional.
- Marginally $p(\lambda|X)$ does not concentrate.
- Marginally, $p(\gamma|X)$ concentrates.

Example: Exponential families with random effects (REs) λ and fixed effects γ .

If the observations per random effect remains bounded as $N \rightarrow \infty$, then

- Parameter λ (“local”) grows in dimension with N .
- Parameter γ (“global”) is finite-dimensional.
- Marginally $p(\lambda|X)$ does not concentrate.
- Marginally, $p(\gamma|X)$ concentrates.

In general, we cannot hope for an asymptotic analysis of $\mathbb{E}_{p(\lambda, \gamma|X)} [f(\lambda)]$.

High dimensional problems

Example: Exponential families with random effects (REs) λ and fixed effects γ .

If the observations per random effect remains bounded as $N \rightarrow \infty$, then

- Parameter λ (“local”) grows in dimension with N .
- Parameter γ (“global”) is finite-dimensional.
- Marginally $p(\lambda|X)$ does not concentrate.
- Marginally, $p(\gamma|X)$ concentrates.

In general, we cannot hope for an asymptotic analysis of $\mathbb{E}_{p(\lambda, \gamma|X)} [f(\lambda)]$.

Can we save the approximation when *some* parameters concentrate?

Does the residual vanish asymptotically for $w_n \mapsto \mathbb{E}_{p(\gamma|X, w_n)} [\gamma]$?

High dimensional problems

We assume that $p(\gamma|X)$ concentrates but $p(\lambda|X)$ does not. By our series expansion:

$$\begin{aligned} & \mathbb{E}_{p(\gamma, \lambda|X, w_n)} [\gamma] - \mathbb{E}_{p(\gamma, \lambda|X)} [\gamma] = \\ & \quad \psi_n(w_n - 1) + \mathcal{E}(w_n) \\ = & \mathbb{E}_{p(\gamma, \lambda|X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma, \lambda|X, \tilde{w}_n)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)^2] (w_n - 1)^2 \end{aligned}$$

$$\psi_n = O_p(N^{-1})$$

$$\mathcal{E}(w_n) = O_p(N^{-1})$$

High dimensional problems

We assume that $p(\gamma|X)$ concentrates but $p(\lambda|X)$ does not. By our series expansion:

$$\begin{aligned}
 & \mathbb{E}_{p(\gamma, \lambda|X, w_n)} [\gamma] - \mathbb{E}_{p(\gamma, \lambda|X)} [\gamma] = \\
 & \quad \psi_n(w_n - 1) + \mathcal{E}(w_n) \\
 & = \mathbb{E}_{p(\gamma, \lambda|X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma, \lambda|X, \tilde{w}_n)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)^2] (w_n - 1)^2 \\
 & = \mathbb{E}_{p(\gamma|X)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|\gamma, X)} [\bar{\ell}_n(\gamma, \lambda)]}_{F_1(\gamma)} \right] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma|X, \tilde{w}_n)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|X, \gamma, \tilde{w}_n)} [\bar{\ell}_n(\gamma, \lambda)^2]}_{F_2(\gamma)} \right] (w_n - 1)
 \end{aligned}$$

$$\psi_n = O_p(N^{-1})$$

$$\mathcal{E}(w_n) = O_p(N^{-1})$$

High dimensional problems

We assume that $p(\gamma|X)$ concentrates but $p(\lambda|X)$ does not. By our series expansion:

$$\begin{aligned}
 & p(\gamma, \lambda|X, w_n) \mathbb{E}[\gamma] - p(\gamma, \lambda|X) \mathbb{E}[\gamma] = \\
 & \quad \psi_n(w_n - 1) + \mathcal{E}(w_n) \\
 = & \mathbb{E}_{p(\gamma, \lambda|X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma, \lambda|X, \tilde{w}_n)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)^2] (w_n - 1)^2 \\
 = & \mathbb{E}_{p(\gamma|X)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|\gamma, X)} [\bar{\ell}_n(\gamma, \lambda)]}_{F_1(\gamma)} \right] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma|X, \tilde{w}_n)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|X, \gamma, \tilde{w}_n)} [\bar{\ell}_n(\gamma, \lambda)^2]}_{F_2(\gamma)} \right] (w_n - 1) \\
 = & \underbrace{\mathbb{E}_{p(\gamma|X)} [\bar{\gamma} F_1(\gamma)]}_{O_p(N^{-1})} (w_n - 1) + \frac{1}{2} \underbrace{\mathbb{E}_{p(\gamma|X, \tilde{w}_n)} [\bar{\gamma} F_2(\gamma)]}_{O_p(N^{-1})} (w_n - 1)^2 \\
 & \text{(by } p(\gamma|X) \text{ concentration)} \quad \text{(by } p(\gamma|X) \text{ concentration)} \\
 \Rightarrow & \psi_n = O_p(N^{-1}) \quad \mathcal{E}(w_n) = O_p(N^{-1})
 \end{aligned}$$

High dimensional problems

We assume that $p(\gamma|X)$ concentrates but $p(\lambda|X)$ does not. By our series expansion:

$$\begin{aligned}
 & p(\gamma, \lambda|X, w_n) [\gamma] - p(\gamma, \lambda|X) [\gamma] = \\
 & \quad \psi_n(w_n - 1) + \mathcal{E}(w_n) \\
 = & \quad \mathbb{E}_{p(\gamma, \lambda|X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma, \lambda|X, \tilde{w}_n)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)^2] (w_n - 1)^2 \\
 = & \quad \mathbb{E}_{p(\gamma|X)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|\gamma, X)} [\bar{\ell}_n(\gamma, \lambda)]}_{F_1(\gamma)} \right] (w_n - 1) + \frac{1}{2} \mathbb{E}_{p(\gamma|X, \tilde{w}_n)} \left[\bar{\gamma} \underbrace{\mathbb{E}_{p(\lambda|X, \gamma, \tilde{w}_n)} [\bar{\ell}_n(\gamma, \lambda)^2]}_{F_2(\gamma)} \right] (w_n - 1) \\
 = & \quad \underbrace{\mathbb{E}_{p(\gamma|X)} [\bar{\gamma} F_1(\gamma)]}_{O_p(N^{-1})} (w_n - 1) + \frac{1}{2} \underbrace{\mathbb{E}_{p(\gamma|X, \tilde{w}_n)} [\bar{\gamma} F_2(\gamma)]}_{O_p(N^{-1})} (w_n - 1)^2 \\
 & \quad \text{(by } p(\gamma|X) \text{ concentration)} \quad \text{(by } p(\gamma|X) \text{ concentration)} \\
 \Rightarrow & \quad \psi_n = O_p(N^{-1}) \quad \mathcal{E}(w_n) = O_p(N^{-1})
 \end{aligned}$$

Corollary [Giordano and Broderick, 2023]:

In general, $w_n \mapsto N \left(\mathbb{E}_{p(\gamma|X, w_n)} [\gamma] - \mathbb{E}_{p(\gamma|X)} [\gamma] \right)$ remains non-linear as $N \rightarrow \infty$.

Example: Poisson regression with Gamma-distributed random effects

For $g = 1, \dots, G$, $\lambda_g \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$ for fixed α, β

For $n = 1, \dots, N$, $g_n \stackrel{iid}{\sim} \text{Categorical}(1, \dots, G)$, $y_n | \lambda_n, \gamma, g_n \stackrel{iid}{\sim} \text{Poisson}(\gamma \lambda_{g_n})$.

$x_n = (y_n, g_n)$ are IID given λ, γ . Write $\log p(X | \lambda, \gamma, w) = \sum_{n=1}^N w_n \ell_n(\lambda, \gamma)$.

Experiments

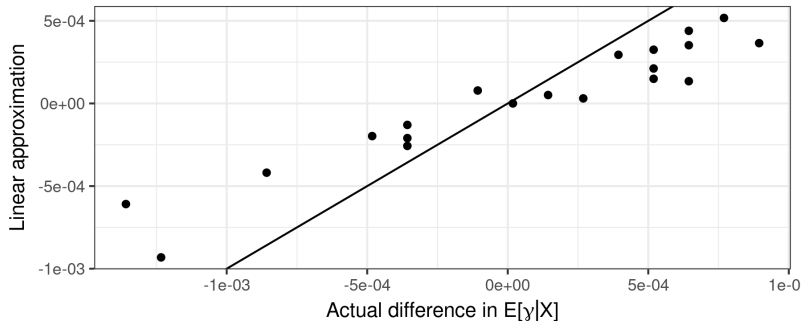
Example: Poisson regression with Gamma-distributed random effects

For $g = 1, \dots, G$, $\lambda_g \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$ for fixed α, β

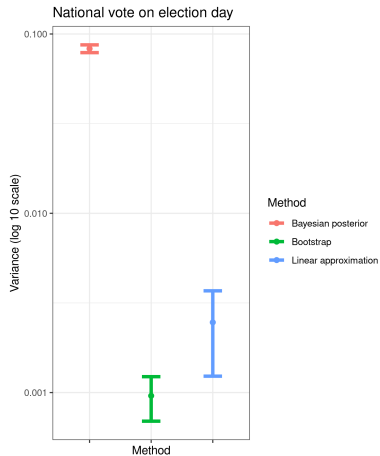
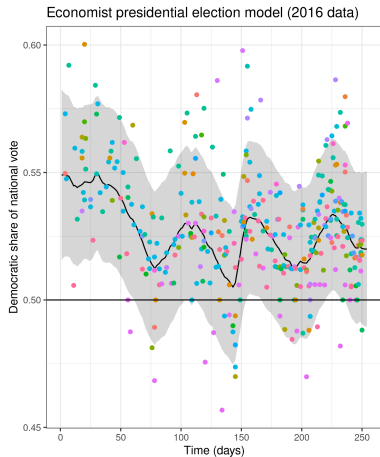
For $n = 1, \dots, N$, $g_n \stackrel{iid}{\sim} \text{Categorical}(1, \dots, G)$, $y_n | \lambda_n, \gamma, g_n \stackrel{iid}{\sim} \text{Poisson}(\gamma \lambda_{g_n})$.

$x_n = (y_n, g_n)$ are IID given λ, γ . Write $\log p(X | \lambda, \gamma, w) = \sum_{n=1}^N w_n \ell_n(\lambda, \gamma)$.

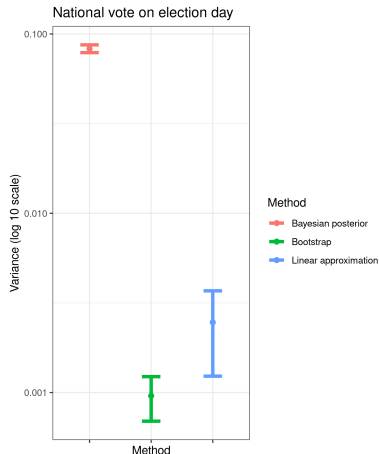
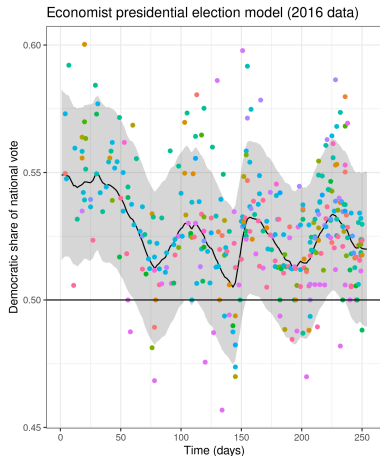
Poisson random effect model
leaving out single datapoints with $N = 800$



Observations and consequences



Observations and consequences



- We often use models of the form $p(\gamma, \lambda|X)$.
- Even if the error $\mathcal{E}(w)$ does not vanish, it can still be small enough in practice.
... Especially given the linear approximation's huge computational advantage.

Preprint: Giordano and Broderick [2023] (arXiv:2305.06466)

(The preprint focuses on variance estimation, the present results are found in the proofs.)

- A. Gelman and M. Heidemanns. The Economist: Forecasting the US elections., 2020. URL <https://projects.economist.com/us-2020-forecast/president>. Data and model accessed Oct., 2020.
- R. Giordano and T. Broderick. The Bayesian infinitesimal jackknife for variance. *arXiv preprint arXiv:2305.06466*, 2023.
- J. Huggins and J. Miller. Reproducible model selection using bagged posteriors. *Bayesian Analysis*, 18(1):79–104, 2023.
- R. Kass, L. Tierney, and J. Kadane. The validity of posterior expansions based on Laplace’s method. *Bayesian and Likelihood Methods in Statistics and Econometrics*, 1990.
- B. Kleijn and A. Van der Vaart. The Bernstein-von-Mises theorem under misspecification. *Electronic Journal of Statistics*, 6: 354–381, 2012.

How can we use the approximation?

How can we use the approximation?

Cross validation. Let $w_{(-n)}$ leave out point n , and loss $f(\theta) = -\ell(x_n|\theta)$.

$$\text{LOO CV loss at point } n = \mathbb{E}_{p(\theta|x, w_{(-n)})} [f(\theta)] \approx \mathbb{E}_{p(\theta|x)} [f(\theta)] - \psi_n$$

How can we use the approximation?

Cross validation. Let $w_{(-n)}$ leave out point n , and loss $f(\theta) = -\ell(x_n|\theta)$.

$$\text{LOO CV loss at point } n = \mathbb{E}_{p(\theta|x, w_{(-n)})} [f(\theta)] \approx \mathbb{E}_{p(\theta|x)} [f(\theta)] - \psi_n$$

Example: Approximate bootstrap.

Draw bootstrap weights $w \sim p(w) = \text{Multinomial}(N, N^{-1})$.

$$\begin{aligned} \text{Bootstrap variance} &= \text{Var}_{p(w)} \left(\mathbb{E}_{p(\theta|x, w)} [f(\theta)] \right) \\ &\approx \text{Var}_{p(w)} \left(\mathbb{E}_{p(\theta|x)} [f(\theta)] + \psi_n (w_n - 1) \right) \\ &= \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2. \end{aligned}$$

How can we use the approximation?

Cross validation. Let $w_{(-n)}$ leave out point n , and loss $f(\theta) = -\ell(x_n|\theta)$.

$$\text{LOO CV loss at point } n = \mathbb{E}_{p(\theta|x, w_{(-n)})} [f(\theta)] \approx \mathbb{E}_{p(\theta|x)} [f(\theta)] - \psi_n$$

Example: Approximate bootstrap.

Draw bootstrap weights $w \sim p(w) = \text{Multinomial}(N, N^{-1})$.

$$\begin{aligned} \text{Bootstrap variance} &= \text{Var}_{p(w)} \left(\mathbb{E}_{p(\theta|x, w)} [f(\theta)] \right) \\ &\approx \text{Var}_{p(w)} \left(\mathbb{E}_{p(\theta|x)} [f(\theta)] + \psi_n(w_n - 1) \right) \\ &= \sum_{n=1}^N \left(\psi_n - \bar{\psi} \right)^2. \end{aligned}$$

Influential subsets: Approximate maximum influence perturbation (AMIP).

Let $W_{(-K)}$ denote weights leaving out K points.

$$\max_{w \in W_{(-K)}} \left(\mathbb{E}_{p(\theta|x, w)} [f(\theta)] - \mathbb{E}_{p(\theta|x)} [f(\theta)] \right) \approx - \sum_{n=1}^K \psi_{(n)}.$$

Example: A negative binomial model

Consider $p(X|\gamma) = \prod_{n=1}^N \text{NegativeBinomial}(x_n|\gamma)$. Here, $\theta = \gamma$ is a scalar.

Example: A negative binomial model

Consider $p(X|\gamma) = \prod_{n=1}^N \text{NegativeBinomial}(x_n|\gamma)$. Here, $\theta = \gamma$ is a scalar.

As $N \rightarrow \infty$, $p(\gamma|X)$ concentrates at rate $1/\sqrt{N}$ (Bernstein–von Mises).

$$\Rightarrow N \left(\mathbb{E}_{p(\gamma|X, w_n)} [\gamma] - \mathbb{E}_{p(\gamma|X)} [\gamma] \right) = \psi_n(w_n - 1) + O_p(N^{-1}).$$

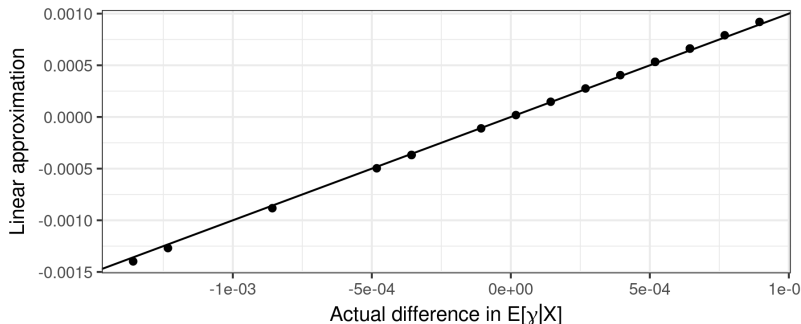
Example: A negative binomial model

Consider $p(X|\gamma) = \prod_{n=1}^N \text{NegativeBinomial}(x_n|\gamma)$. Here, $\theta = \gamma$ is a scalar.

As $N \rightarrow \infty$, $p(\gamma|X)$ concentrates at rate $1/\sqrt{N}$ (Bernstein–von Mises).

$$\Rightarrow N \left(\mathbb{E}_{p(\gamma|X, w_n)}[\gamma] - \mathbb{E}_{p(\gamma|X)}[\gamma] \right) = \psi_n(w_n - 1) + O_p(N^{-1}).$$

Negative Binomial model
leaving out single datapoints with $N = 800$



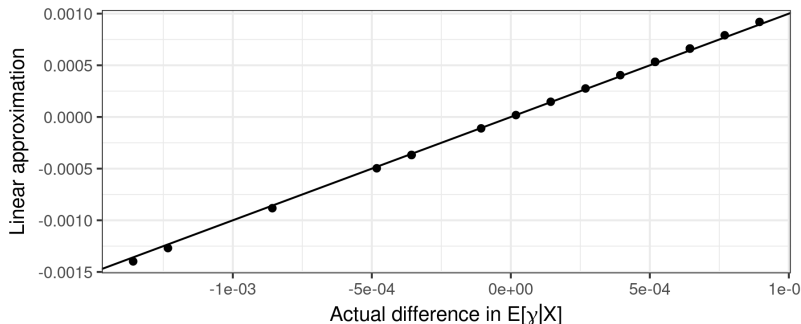
Example: A negative binomial model

Consider $p(X|\gamma) = \prod_{n=1}^N \text{NegativeBinomial}(x_n|\gamma)$. Here, $\theta = \gamma$ is a scalar.

As $N \rightarrow \infty$, $p(\gamma|X)$ concentrates at rate $1/\sqrt{N}$ (Bernstein–von Mises).

$$\Rightarrow N \left(\mathbb{E}_{p(\gamma|X, w_n)}[\gamma] - \mathbb{E}_{p(\gamma|X)}[\gamma] \right) = \psi_n(w_n - 1) + O_p(N^{-1}).$$

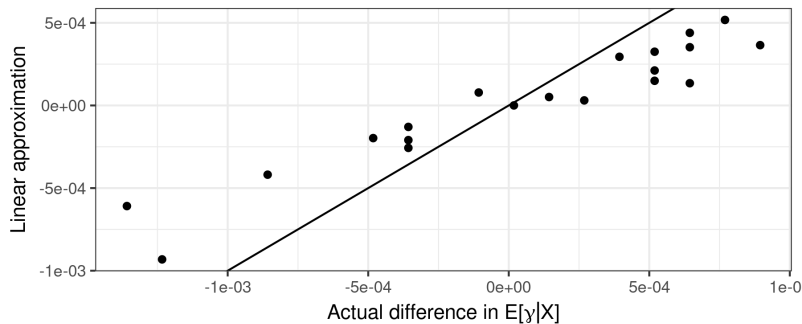
Negative Binomial model
leaving out single datapoints with $N = 800$



Problem: Most computationally hard Bayesian problems don't concentrate.

Example: **Poisson model with random effects (REs) λ and fixed effect γ .**

Poisson random effect model
leaving out single datapoints with $N = 800$



A contradiction?

Negative binomial observations.

Asymptotically linear in w .

Poisson observations with random effects.

Asymptotically non-linear in w .

A contradiction?

Negative binomial observations.

Asymptotically linear in w .

Poisson observations with random effects.

Asymptotically non-linear in w .

With a constant regressor, Gamma REs, and one RE per observation,
these are the same model, with the same $p(\gamma|X)$.

Is $\mathbb{E}_{p(\gamma|X,w)} [\gamma]$ linear in the data weights or not?

A contradiction?

Negative binomial observations.

Asymptotically linear in w .

$$\log p(X|\gamma, w^m) = \sum_{n=1}^N w_n^m \log p(x_n|\gamma)$$

Poisson observations with random effects.

Asymptotically non-linear in w .

$$\log p(X|\gamma, \lambda, w^c) = \sum_{n=1}^N w_n^c \log p(x_n|\lambda, \gamma)$$

With a constant regressor, Gamma REs, and one RE per observation, these are the same model, with the same $p(\gamma|X)$.

Is $\mathbb{E}_{p(\gamma|X, w)} [\gamma]$ **linear in the data weights** or not?

Trick question! We weight a log likelihood contribution, not a datapoint.

The two weightings are not equivalent in general.

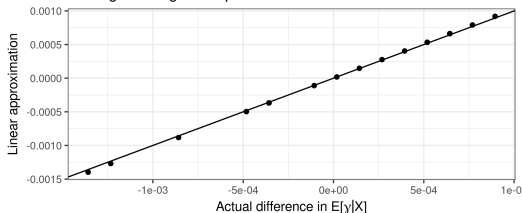
Experimental results

Our results were actually computed on **identical datasets** with $G = N$ and $g_n = n$.

Uses $\log p(x_n | \gamma)$:

$$\psi_n = \mathbb{E}_{p(\gamma|X)} [\bar{\gamma} \bar{\ell}_n(\gamma)]$$

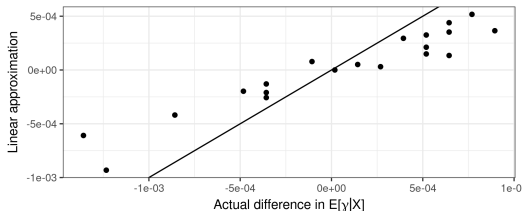
Negative Binomial model
leaving out single datapoints with $N = 800$



Uses $\log p(x_n | \gamma, \lambda)$:

$$\psi_n = \mathbb{E}_{p(\gamma, \lambda|X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)]$$

Poisson random effect model
leaving out single datapoints with $N = 800$



Experimental results

Our results were actually computed on **identical datasets** with $G = N$ and $g_n = n$.

Uses $\log p(x_n | \gamma)$:

$$\psi_n = \mathbb{E}_{p(\gamma|X)} [\bar{\gamma} \bar{\ell}_n(\gamma)]$$

Not computable from

$$\gamma, \lambda \sim p(\gamma, \lambda | X)$$

in general.

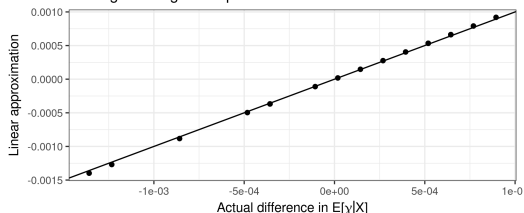
Uses $\log p(x_n | \gamma, \lambda)$:

$$\psi_n = \mathbb{E}_{p(\gamma, \lambda | X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)]$$

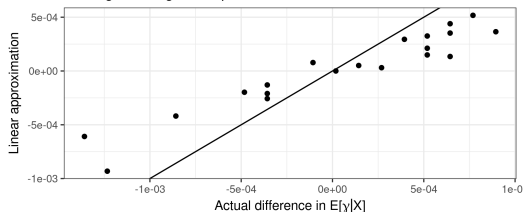
Computable from

$$\gamma, \lambda \sim p(\gamma, \lambda | X).$$

Negative Binomial model
leaving out single datapoints with $N = 800$



Poisson random effect model
leaving out single datapoints with $N = 800$



Experimental results

Our results were actually computed on **identical datasets** with $G = N$ and $g_n = n$.

Uses $\log p(x_n | \gamma)$:

$$\psi_n = \mathbb{E}_{p(\gamma|X)} [\bar{\gamma} \bar{\ell}_n(\gamma)]$$

Not computable from

$$\gamma, \lambda \sim p(\gamma, \lambda | X)$$

in general.

Uses $\log p(x_n | \gamma, \lambda)$:

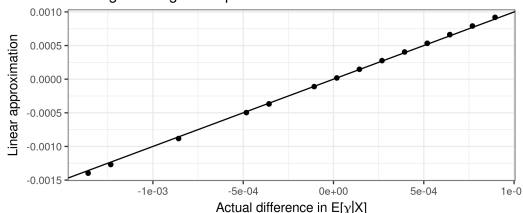
$$\psi_n = \mathbb{E}_{p(\gamma, \lambda | X)} [\bar{\gamma} \bar{\ell}_n(\gamma, \lambda)]$$

Computable from

$$\gamma, \lambda \sim p(\gamma, \lambda | X).$$

May still be useful when $p(\lambda | X)$ is *somewhat* concentrated.

Negative Binomial model
leaving out single datapoints with $N = 800$



Poisson random effect model
leaving out single datapoints with $N = 800$

